

# CompSci 516

# Database Systems

## Lecture 5

# Relational Algebra And Normalization

Guest Lecture: Junyang Gao

# Announcements

- Lab-2 on RA on Thursday
  - Do not forget your laptop!
- Homework-1 (Part-1 and 2) have been posted on sakai
  - First deadline next Tuesday: 09/17
  - Parsing XML will take time!
- Next week:
  - Revisit Relational Calculus!
  - New topic: Database internals and indexes!

# Today's topic

- Relational Algebra
- Normalization

Acknowledgement:

The following slides have been created adapting the instructor material of the [RG] book provided by the authors

Dr. Ramakrishnan and Dr. Gehrke.

# A Quick Recap

Likes(drinker, beer)  
Frequents(drinker, bar)  
Serves(bar, beer)

- The “Drinker-Beer-Bar” example

**Query:** Find drinkers that like some beer (so much) that they frequent all bars that serve it

$$Q(x) = \exists y. \text{Likes}(x, y) \wedge \forall z. (\text{Serves}(z, y) \Rightarrow \text{Frequents}(x, z))$$

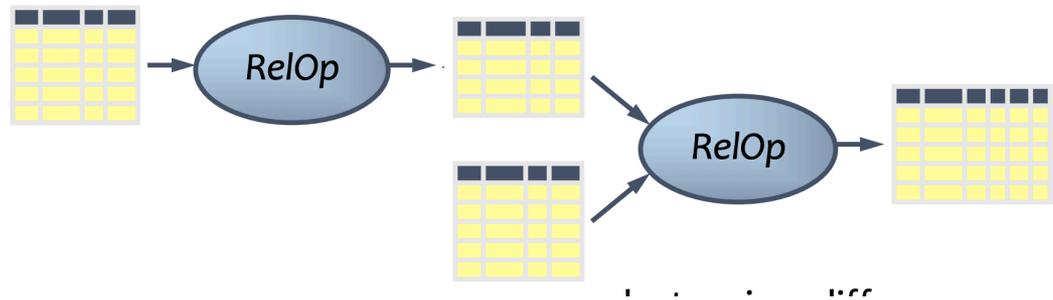
```
SELECT DISTINCT L.drinker
FROM Likes L
WHERE not exists
  (SELECT S.bar
   FROM Serves S
   WHERE L.beer=S.beer
    AND not exists (SELECT *
                   FROM Frequents F
                   WHERE F.drinker=L.drinker
                   AND F.bar=S.bar))
```

# Relational Algebra (RA)

# Relational Algebra

- A language for querying relational data based on “operators”
- Takes one or more relations as input, and produces a relation as output

- operator
- operand
- semantic
- so an algebra!



- Since each operation returns a relation, **operations can be composed**
  - Algebra is “closed”

# Relational Algebra

- Basic operations:
  - Selection ( $\sigma$ ) Selects a subset of rows from relation
  - Projection ( $\pi$ ) Deletes unwanted columns from relation.
  - Cross-product ( $\times$ ) Allows us to combine two relations.
  - Set-difference ( $-$ ) Tuples in reln. 1, but not in reln. 2.
  - Union ( $\cup$ ) Tuples in reln. 1 or in reln. 2.
- Additional operations:
  - Intersection ( $\cap$ )
  - join  $\bowtie$
  - division( $/$ )
  - renaming ( $\rho$ )
  - Not essential, but (very) useful, especially join!

# Example Schema and Instances

Sailors(sid, sname, rating, age)

Boats(bid, bname, color)

Reserves(sid, bid, day)

*S1*

<u>sid</u>	sname	rating	age
22	dustin	7	45.0
31	lubber	8	55.5
58	rusty	10	35.0

*S2*

<u>sid</u>	sname	rating	age
28	yuppy	9	35.0
31	lubber	8	55.5
44	guppy	5	35.0
58	rusty	10	35.0

*R1*

<u>sid</u>	<u>bid</u>	<u>day</u>
22	101	10/10/96
58	103	11/12/96

# Projection

S2

<u>sid</u>	sname	rating	age
28	yuppy	9	35.0
31	lubber	8	55.5
44	guppy	5	35.0
58	rusty	10	35.0

- Deletes attributes that are not in projection list.
- Schema** of result contains exactly the fields in the projection list, with the same names that they had in the (single) input relation.

sname	rating
yuppy	9
lubber	8
guppy	5
rusty	10

$\pi_{sname, rating}(S2)$

age
35.0
55.5
35.0
35.0



$\pi_{age}(S2)$

# Projection

S2

<u>sid</u>	sname	rating	age
28	yuppy	9	35.0
31	lubber	8	55.5
44	guppy	5	35.0
58	rusty	10	35.0

- Deletes attributes that are not in projection list.
- Schema** of result contains exactly the fields in the projection list, with the same names that they had in the (single) input relation.
- Projection operator has to **eliminate duplicates** (Set semantic!)
  - Note: real systems typically don't do duplicate elimination unless the user explicitly asks for it (performance)

sname	rating
yuppy	9
lubber	8
guppy	5
rusty	10

$\pi_{sname, rating}(S2)$

age
35.0
55.5

$\pi_{age}(S2)$

# Selection

- Selects rows that satisfy **selection condition**
- No duplicates in result
  - Because input is a set!
- Schema of result identical to schema of (single) input relation

S2

<u>sid</u>	sname	rating	age
28	yuppy	9	35.0
31	lubber	8	55.5
44	guppy	5	35.0
58	rusty	10	35.0

sid	sname	rating	age
28	yuppy	9	35.0
58	rusty	10	35.0

$$\sigma_{rating > 8}(S2)$$

# Composition of Operators

- Result relation can be the input for another relational algebra operation
  - Operator composition

S2

<u>sid</u>	sname	rating	age
28	yuppy	9	35.0
31	lubber	8	55.5
44	guppy	5	35.0
58	rusty	10	35.0

What do we get here?

$$\pi_{sname, rating}(\sigma_{rating > 8}(S2))$$

# Composition of Operators

- Result relation can be the input for another relational algebra operation
  - Operator composition

S2

<u>sid</u>	sname	rating	age
28	yuppy	9	35.0
31	lubber	8	55.5
44	guppy	5	35.0
58	rusty	10	35.0

sname	rating
yuppy	9
rusty	10

What do we get here?

$$\pi_{sname, rating}(\sigma_{rating > 8}(S2))$$

# Union, Intersection, Set-Difference

S1

<u>sid</u>	sname	rating	age
22	dustin	7	45.0
31	lubber	8	55.5
58	rusty	10	35.0

S2

<u>sid</u>	sname	rating	age
28	yuppy	9	35.0
31	lubber	8	55.5
44	guppy	5	35.0
58	rusty	10	35.0

- All of these operations take two input relations, which must be **union-compatible**:
  - Same number of fields.
  - “Corresponding” fields must have the same type and same schema as the inputs

<u>sid</u>	sname	rating	age
22	dustin	7	45.0
31	lubber	8	55.5
58	rusty	10	35.0
44	guppy	5	35.0
28	yuppy	9	35.0

You would lose points if your relations in  $\cup, \cap, -$  are not union compatible!

# Union, Intersection, Set-Difference

S1

<u>sid</u>	sname	rating	age
22	dustin	7	45.0
31	lubber	8	55.5
58	rusty	10	35.0

S2

<u>sid</u>	sname	rating	age
28	yuppy	9	35.0
31	lubber	8	55.5
44	guppy	5	35.0
58	rusty	10	35.0

- Note: no duplicate
  - “Set semantic”
  - SQL: **UNION**
  - But SQL allows “bag semantic” as well:  
**UNION ALL**

<u>sid</u>	sname	rating	age
22	dustin	7	45.0
31	lubber	8	55.5
58	rusty	10	35.0
44	guppy	5	35.0
28	yuppy	9	35.0

# Union, Intersection, Set-Difference

*S1*

<u>sid</u>	sname	rating	age
22	dustin	7	45.0
31	lubber	8	55.5
58	rusty	10	35.0

*S2*

<u>sid</u>	sname	rating	age
28	yuppy	9	35.0
31	lubber	8	55.5
44	guppy	5	35.0
58	rusty	10	35.0

<u>sid</u>	sname	rating	age
22	dustin	7	45.0

$S1 - S2$

<u>sid</u>	sname	rating	age
31	lubber	8	55.5
58	rusty	10	35.0

$S1 \cap S2$

# Cross-Product

- Each row of S1 is paired with each row of R.
- **Result schema** has one field per field of S1 and R, with field names 'inherited' if possible.
  - Conflict: Both S1 and R have a field called sid.

<u>sid</u>	sname	rating	age
22	dustin	7	45.0
31	lubber	8	55.5
58	rusty	10	35.0

<u>sid</u>	<u>bid</u>	<u>day</u>
22	101	10/10/96
58	103	11/12/96

(sid)	sname	rating	age	(sid)	bid	day
22	dustin	7	45.0	22	101	10/10/96
22	dustin	7	45.0	58	103	11/12/96
31	lubber	8	55.5	22	101	10/10/96
31	lubber	8	55.5	58	103	11/12/96
58	rusty	10	35.0	22	101	10/10/96
58	rusty	10	35.0	58	103	11/12/96

# Renaming Operator $\rho$

$$(\rho_{\text{sid} \rightarrow \text{sid1}} S1) \times (\rho_{\text{sid} \rightarrow \text{sid2}} R1)$$

Different syntaxes are used  
You can use any of these

or

$$\rho(C(1 \rightarrow \text{sid1}, 5 \rightarrow \text{sid2}), S1 \times R1)$$

C is the  
new relation  
name

<del>sid1</del> (sid)	sname	rating	age	<del>sid2</del> (sid)	bid	day
22	dustin	7	45.0	22	101	10/10/96
22	dustin	7	45.0	58	103	11/12/96
31	lubber	8	55.5	22	101	10/10/96
31	lubber	8	55.5	58	103	11/12/96
58	rusty	10	35.0	22	101	10/10/96
58	rusty	10	35.0	58	103	11/12/96

- In general, can use  $\rho(\langle \text{Temp} \rangle, \langle \text{RA-expression} \rangle)$

# Joins

$$R \bowtie_c S = \sigma_c (R \times S)$$

(sid)	sname	rating	age	(sid)	bid	day
22	dustin	7	45.0	58	103	11/12/96
31	lubber	8	55.5	58	103	11/12/96

$$S1 \bowtie_{S1.sid < R1.sid} R1$$

- Result schema same as that of cross-product.
- Fewer tuples than cross-product, might be able to compute more efficiently
  - We will do join algorithms later

# Find names of sailors who've reserved boat #103

Sailors(sid, sname, rating, age)

Boats(bid, bname, color)

Reserves(sid, bid, day)

# Find names of sailors who've reserved boat #103

Sailors(sid, sname, rating, age)

Boats(bid, bname, color)

Reserves(sid, bid, day)

# Find names of sailors who've reserved boat #103

Sailors(sid, sname, rating, age)

Boats(bid, bname, color)

Reserves(sid, bid, day)

- **Solution 1:**  $\pi_{sname}((\sigma_{bid=103} Reserves) \bowtie Sailors)$
- **Solution 2:**  $\pi_{sname}(\sigma_{bid=103}(Reserves \bowtie Sailors))$

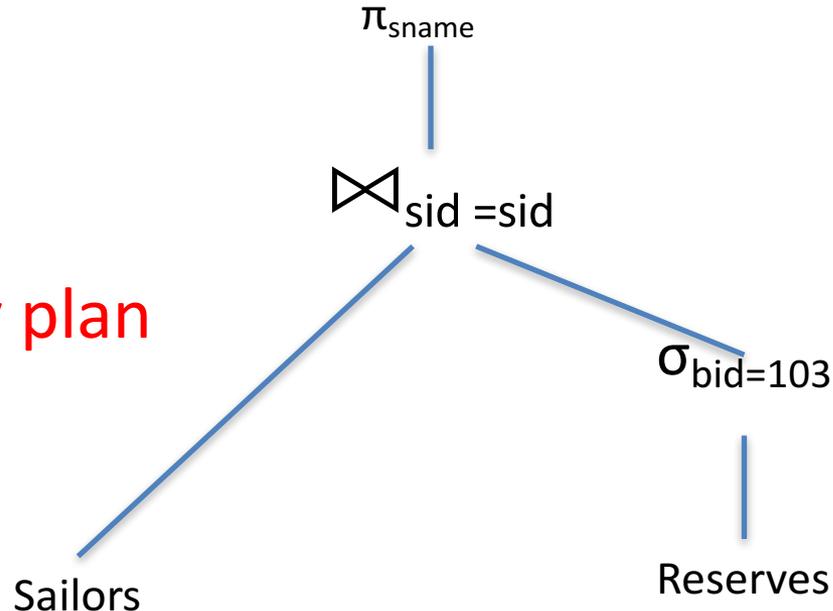
# Expressing an RA expression as a Tree

Sailors(sid, sname, rating, age)

Boats(bid, bname, color)

Reserves(sid, bid, day)

Also called a  
**logical query plan**



You can use either  
RA expression or a logical  
query plan in exams

But in the lab on Thursday  
you will use RA expressions  
In RADB and RAtest tools

$$\pi_{sname}((\sigma_{bid=103} \text{Reserves}) \bowtie \text{Sailors})$$

# Find sailors who've reserved a red or a green boat

Sailors(sid, sname, rating, age)

Boats(bid, bname, color)

Reserves(sid, bid, day)

Use of rename operation

# Find sailors who've reserved a red or a green boat

Sailors(sid, sname, rating, age)

Boats(bid, bname, color)

Reserves(sid, bid, day)

Use of rename operation

- Can identify all red or green boats, then find sailors who've reserved one of these boats:

$$\rho (Tempboats, (\sigma_{color='red' \vee color='green'} Boats))$$
$$\pi_{sname}(Tempboats \bowtie Reserves \bowtie Sailors)$$

Can also define Tempboats using union

Try the “AND” version yourself

# What about aggregates?

Sailors(sid, sname, rating, age)

Boats(bid, bname, color)

Reserves(sid, bid, day)

- Supported by extended relational algebra
- $\gamma_{age, avg(rating)} \rightarrow avgr$  Sailors
- Also extended to “bag semantic”: allow duplicates
  - Take into account cardinality
  - R and S have tuple t resp. m and n times
  - $R \cup S$  has t m+n times
  - $R \cap S$  has t min(m, n) times
  - $R - S$  has t max(0, m-n) times
  - sorting( $\tau$ ), duplicate removal ( $\delta$ ) operators

# Example: Aggregates in RA

Sailors(sid, sname, rating, age)

Boats(bid, bname, color)

Reserves(sid, bid, day)

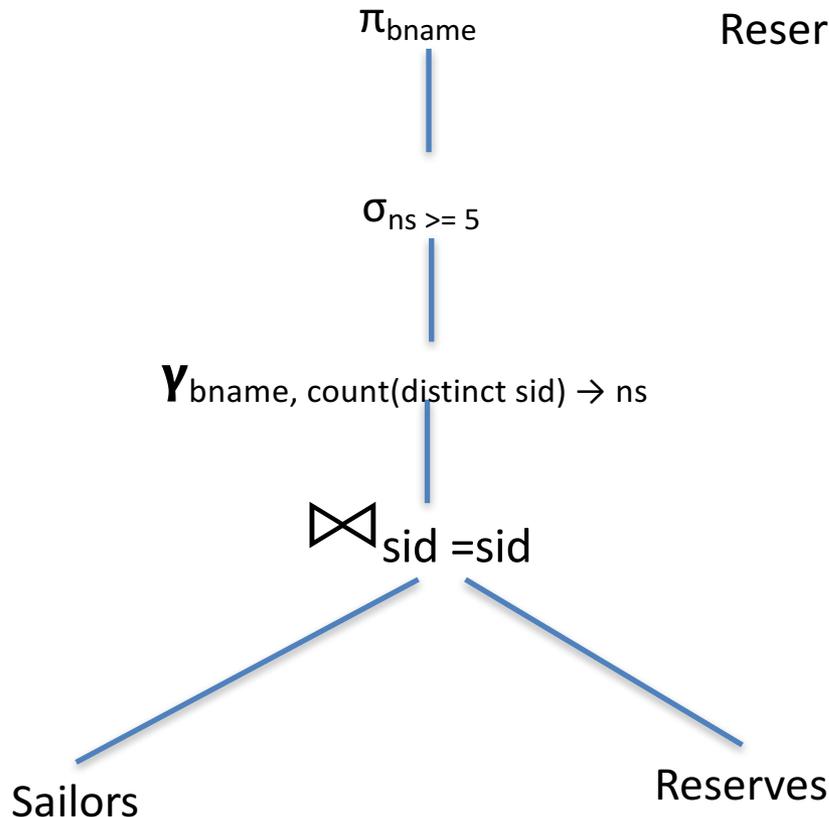
Output bnames that  
have been reserved by at  
least 5 sailors

# Example: Aggregates in RA

Sailors(sid, sname, rating, age)

Boats(bid, bname, color)

Reserves(sid, bid, day)



Output bnames that have been reserved by at least 5 sailors

# Database Normalization

# What will we learn?

- What goes wrong if we have redundant info in a database?
- Why and how should you refine a schema?
- Functional Dependencies – a new kind of integrity constraints (IC)
- Normal Forms
- How to obtain those normal forms

# Example

The list of hourly employees in an organization

<u>ssn (S)</u>	name (N)	lot (L)	rating (R)	hourly-wage (W)	hours-worked (H)
111-11-1111	Attishoo	48	8	10	40
222-22-2222	Smiley	22	8	10	30
333-33-3333	Smethurst	35	5	7	30
444-44-4444	Guldu	35	5	7	32
555-55-5555	Madayan	35	8	10	40

- key = SSN
- Suppose for a given rating, there is only one hourly\_wage value
- Redundancy in the table
- Why is redundancy bad?

# Why is redundancy bad? 1/4

The list of hourly employees in an organization

<u>ssn (S)</u>	name (N)	lot (L)	rating (R)	hourly-wage (W)	hours-worked (H)
111-11-1111	Attishoo	48	8	10	40
222-22-2222	Smiley	22	8	10	30
333-33-3333	Smethurst	35	5	7	30
444-44-4444	Guldu	35	5	7	32
555-55-5555	Madayan	35	8	10	40

## 1. Redundant storage:

- Some information is stored repeatedly
- The rating value 8 corresponds to hourly\_wage 10, which is stored three times

# Why is redundancy bad? 2/4

The list of hourly employees in an organization

<u>ssn (S)</u>	name (N)	lot (L)	rating (R)	hourly-wage (W)	hours-worked (H)
111-11-1111	Attishoo	48	8	10 → 9	40
222-22-2222	Smiley	22	8	10	30
333-33-3333	Smethurst	35	5	7	30
444-44-4444	Guldu	35	5	7	32
555-55-5555	Madayan	35	8	10	40

## 2. Update anomalies

- If one copy of data is updated, an inconsistency is created unless all copies are similarly updated
- Suppose you update the hourly\_wage value in the first tuple using UPDATE statement in SQL -- inconsistency

# Why is redundancy bad? 3/4

The list of hourly employees in an organization

<u>ssn (S)</u>	name (N)	lot (L)	rating (R)	hourly-wage (W)	hours-worked (H)
111-11-1111	Attishoo	48	8	10	40
222-22-2222	Smiley	22	8	10	30
333-33-3333	Smethurst	35	5	7	30
444-44-4444	Guldu	35	5	7	32
555-55-5555	Madayan	35	8	10	40

### 3. Insertion anomalies:

- It may not be possible to store certain information unless some other, unrelated info is stored as well
- We cannot insert a tuple for an employee unless we know the hourly wage for the employee's rating value

# Why is redundancy bad? 4/4

The list of hourly employees in an organization

<u>ssn (S)</u>	name (N)	lot (L)	rating (R)	hourly-wage (W)	hours-worked (H)
111-11-1111	Attishoo	48	8	10	40
222-22-2222	Smiley	22	8	10	30
333-33-3333	Smethurst	35	5	7	30
444-44-4444	Guldu	35	5	7	32
555-55-5555	Madayan	35	8	10	40

## 4. Deletion anomalies:

- It may not be possible delete certain information without losing some other information as well
- If we delete all tuples with a given rating value (Attishoo, Smiley, Madayan), we lose the association between that rating value and its hourly\_wage value

# Nulls may or may not help

<u>ssn (S)</u>	name (N)	lot (L)	rating (R)	hourly-wage (W)	hours-worked (H)
111-11-1111	Attishoo	48	8	10	40
222-22-2222	Smiley	22	8	10	30
333-33-3333	Smethurst	35	5	7	30
444-44-4444	Guldu	35	5	7	32
555-55-5555	Madayan	35	8	10	40

- Does not help redundant storage or update anomalies
- May help insertion and deletion anomalies
  - can insert a tuple with null value in the hourly\_wage field
  - but cannot record hourly\_wage for a rating unless there is such an employee (SSN cannot be null) – same for deletion

# Summary: Redundancy

- Redundancy arises when the schema forces an association between attributes that is “not natural”
- We want schemas that do not permit redundancy
  - at least identify schemas that allow redundancy to make an informed decision (e.g. for performance reasons)
- Null value may or may not help
- **Solution?**
  - **Decomposition of schema!**

# Decomposition: Example-1

<u>ssn (S)</u>	name (N)	lot (L)	rating (R)	hourly-wage (W)	hours-worked (H)
111-11-1111	Attishoo	48	8	10	40
222-22-2222	Smiley	22	8	10	30
333-33-3333	Smethurst	35	5	7	30
444-44-4444	Guldu	35	5	7	32
555-55-5555	Madayan	35	8	10	40

# Decomposition: Example-1

<u>ssn (S)</u>	name (N)	lot (L)	rating (R)	hourly-wage (W)	hours-worked (H)
111-11-1111	Attishoo	48	8	10	40
222-22-2222	Smiley	22	8	10	30
333-33-3333	Smethurst	35	5	7	30
444-44-4444	Guldu	35	5	7	32
555-55-5555	Madayan	35	8	10	40

<u>ssn (S)</u>	name (N)	lot (L)	rating (R)	hours-worked (H)
111-11-1111	Attishoo	48	8	40
222-22-2222	Smiley	22	8	30
333-33-3333	Smethurst	35	5	30
444-44-4444	Guldu	35	5	32
555-55-5555	Madayan	35	8	40

<u>rating</u>	hourly_wage
8	10
5	7

# Decomposition – Example-2

(on twitter)

<i>uid</i>	<i>uname</i>	<i>twitterid</i>	<i>gid</i>	<i>fromDate</i>
142	Bart	@BartJSimpson	dps	1987-04-19
123	Milhouse	@MilhouseVan_	gov	1989-12-17
857	Lisa	@lisasimpson	abc	1987-04-19
857	Lisa	@lisasimpson	gov	1988-09-01
456	Ralph	@ralphwiggum	abc	1991-04-25
456	Ralph	@ralphwiggum	gov	1992-09-01
...	...	...	...	...

- User id
- user name
- Twitter id
- Group id
- Joining Date (to a group)
- Both uid and twitterid are keys

<i>uid</i>	<i>uname</i>	<i>twitterid</i>
142	Bart	@BartJSimpson
123	Milhouse	@MilhouseVan_
857	Lisa	@lisasimpson
456	Ralph	@ralphwiggum
...	...	...

<i>uid</i>	<i>gid</i>	<i>fromDate</i>
142	dps	1987-04-19
123	gov	1989-12-17
857	abc	1987-04-19
857	gov	1988-09-01
456	abc	1991-04-25
456	gov	1992-09-01
...	...	...

# Unnecessary decomposition

<i>uid</i>	<i>uname</i>	<i>twitterid</i>
142	Bart	@BartJSimpson
123	Milhouse	@MilhouseVan_
857	Lisa	@lisasimpson
456	Ralph	@ralphwiggum
...	...	...

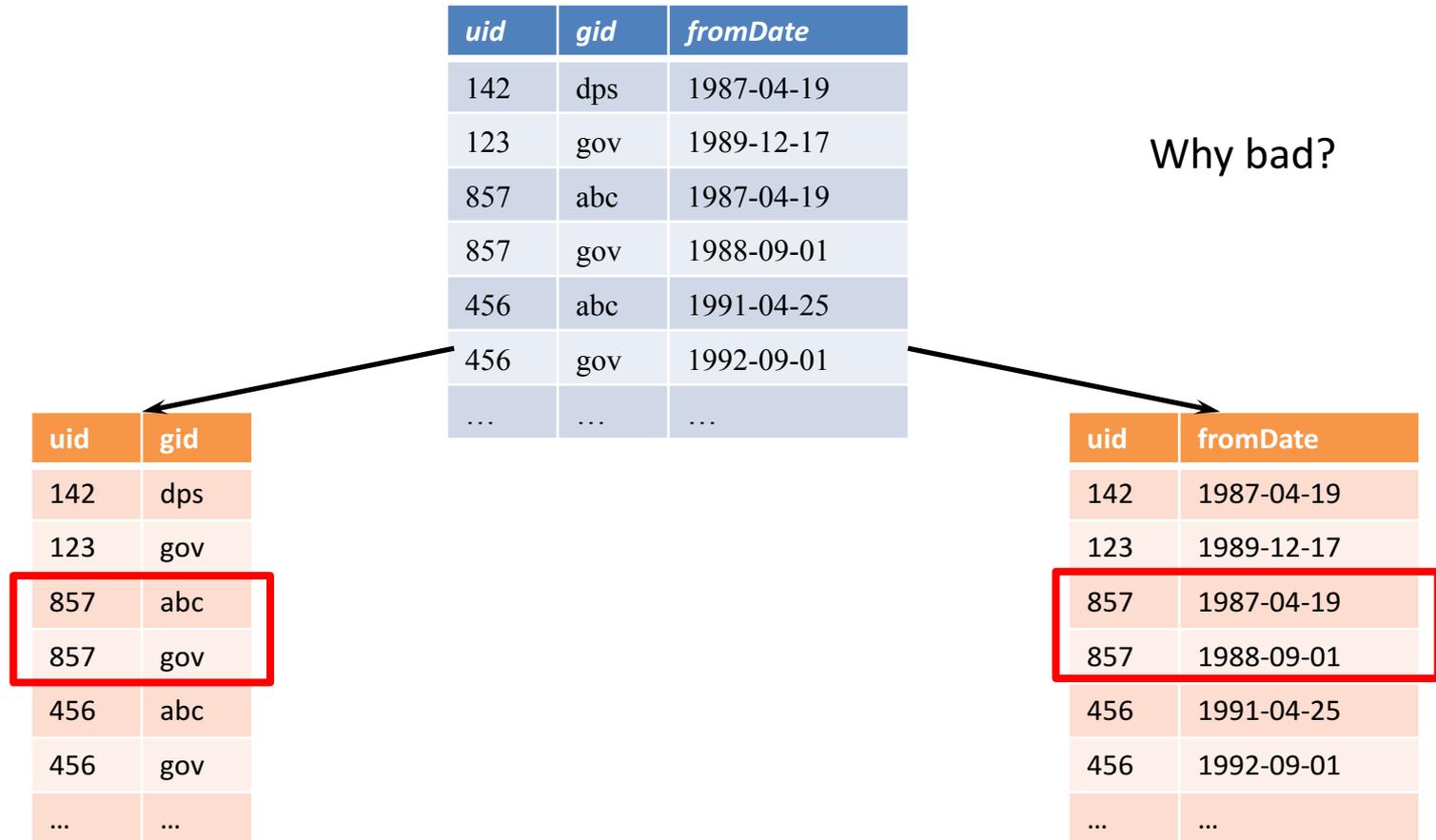
Why?

<i>uid</i>	<i>uname</i>
142	Bart
123	Milhouse
857	Lisa
456	Ralph
...	...

<i>uid</i>	<i>twitterid</i>
142	@BartJSimpson
123	@MilhouseVan_
857	@lisasimpson
456	@ralphwiggum
...	...

- **Still correct:** join returns the original relation
- **Unnecessary:** no redundancy is removed; schema is more complicated (and *uid* is stored twice!)

# Bad decomposition



- Association between *gid* and *fromDate* is lost
- Join returns more rows than the original relation

# Lossless join decomposition

- Decompose relation  $R$  into relations  $S$  and  $T$ 
  - $attrs(R) = attrs(S) \cup attrs(T)$
  - $S = \pi_{attrs(S)}(R)$
  - $T = \pi_{attrs(T)}(R)$
- The decomposition is a **lossless join decomposition** if, given known constraints such as FD's, we can guarantee that  $R = S \bowtie T$
- $R \subseteq S \bowtie T$  or  $R \supseteq S \bowtie T$  ?
- Any decomposition gives  $R \subseteq S \bowtie T$  (why?)
  - A **lossy** decomposition is one with  $R \subset S \bowtie T$

# Loss? But I got more rows!

- “Loss” refers not to the loss of tuples, but to the loss of information
  - Or, the ability to distinguish different original relations

<i>uid</i>	<i>gid</i>	<i>fromDate</i>
142	dps	1987-04-19
123	gov	1989-12-17
857	abc	1988-09-01
857	gov	1987-04-19
456	abc	1991-04-25
456	gov	1992-09-01
...	...	...

No way to tell  
which is the original relation

<i>uid</i>	<i>gid</i>
142	dps
123	gov
857	abc
857	gov
456	abc
456	gov
...	...

<i>uid</i>	<i>fromDate</i>
142	1987-04-19
123	1989-12-17
857	1987-04-19
857	1988-09-01
456	1991-04-25
456	1992-09-01
...	...

# Decompositions should be used judiciously

## 1. Do we need to decompose a relation?

- Several “**normal forms**” exist to identify possible redundancy at different granularity
- If a relation is not in one of them, may need to decompose further

## 2. What are the problems with decomposition?

- **Bad decompositions**: e.g., Lossy decompositions
- **Performance issues** -- decomposition may both
  - help performance (for updates, some queries accessing part of data), or
  - hurt performance (new joins may be needed for some queries)